

Saprothub: Making Protein Modeling Accessible to All Biologists

Jin Su¹, Zhikai Li¹, Chenchen Han¹, Yuyang Zhou¹, Yan He¹, Junjie Shan¹, Xibin Zhou¹, Xing Chang¹, Dacheng Ma², The OPMC⁵, Martin Steinegger³, Sergey Ovchinnikov⁴, Fajie Yuan¹

¹Westlake University.

²Zhejiang Lab.

³Seoul National University.

⁴MIT.

⁵Open Protein Modeling Consortium.

Correspondence to: yuanfajie@westlake.edu.cn;

Abstract

Training and deploying deep learning models pose challenges for users without machine learning (ML) expertise. Saprothub offers a user-friendly platform that democratizes the process of training, utilizing, storing, and sharing protein ML models, fostering collaboration within the biology community—all achievable with just a few clicks, regardless of ML background. At its core, Saproth is an advanced, foundational protein language model. Through its ColabSaproth framework, it supports potentially hundreds of protein training and prediction applications, enabling the co-construction and co-sharing of these trained models. This enhances user engagement and drives community-wide innovation.

Proteins, the building blocks of life, are essential for a myriad of biological processes and pivotal in medical breakthroughs, pharmaceutical innovations, and genetic discoveries [5, 20, 10]. Despite their importance, understanding protein structure and function has been a considerable challenge in the scientific community. The success of AlphaFold2 [22] in CASP14 marks a new era in structural biology by achieving accuracy of protein structure prediction at the experimental margin of error. Simultaneously, large-scale protein language models (PLMs) are driving substantial advances in protein function prediction.

In this context, several notable PLMs have emerged, including the ESM models [36, 25, 54, 59], ProtTrans [6], ProstT5 [13], UniRep [2], Tranception [31], ProGen [26, 30], and ProtGPT2 [7], each demonstrating remarkable efficacy in their respective tasks. However, training and deploying these machine learning (ML) models for proteins often pose major challenges for researchers without ML expertise. These include selecting appropriate model architectures, managing coding intricacies, preprocessing large datasets, training model parameters, evaluating and interpreting results. This complexity often deters non-ML researchers, hindering their active involvement in this domain, especially as AI models continue to grow more sophisticated.

Given these challenges, ColabFold [29] emerged as a pioneering initiative by deploying AlphaFold2 predictions on Google Colaboratory, making protein folding accessible to researchers of all backgrounds. However, ColabFold primarily emphasizes model inference, lacking support for researchers in the complex tasks such as training their own ML models.

In response to this need, we present Saprothub, an easy-to-use platform based on Google Colaboratory tailored for protein function-related predictions. Saprothub empowers researchers by enabling them to create and train their own task-specific models without requiring advanced ML or coding expertise. Additionally, Saprothub distinguishes itself by supporting a broad range of protein function predictions (Fig. 1c, 2b), rather than being confined to a specific task.

Furthermore, we introduce the *Open Protein Modeling Consortium* (OPMC) (Supplementary) with the vision to create a unified repository for decentralized protein prediction models. Within the OPMC framework, researchers can easily share their individually trained models, fostering both the direct use and collaborative construction of protein models. This approach enables continuous learning and improvement of

Partial materials of Saproth was reported in [37], following Nature journal's conference proceedings policy, see [here](#).

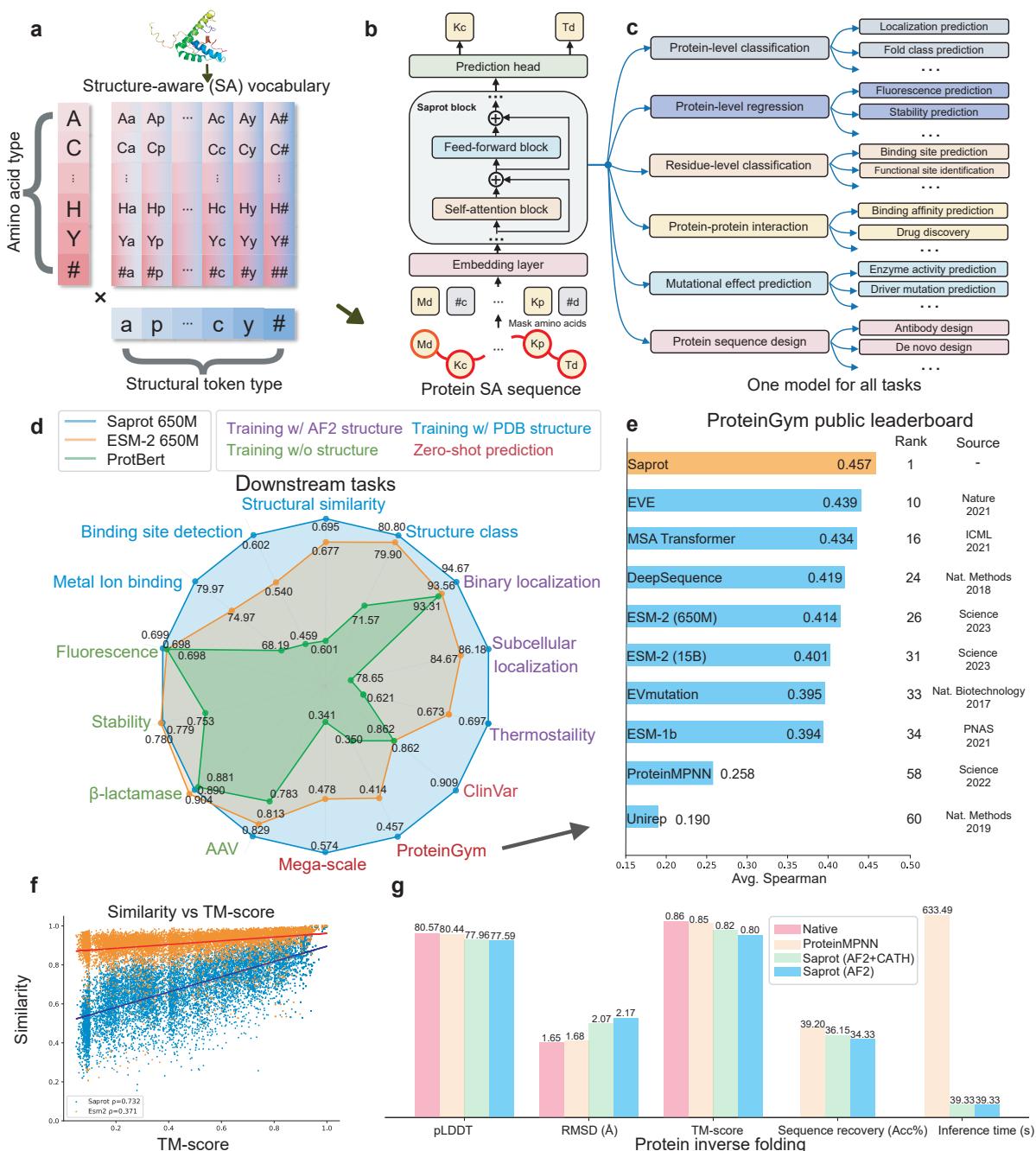


Fig. 1: Illustration of Saprot. **a**, The proposed SA vocabulary. #a, #p,...,#y represent only the 3Di token being visible, while A#, C#,..., Y# represent the AA token being visible. **b**, Architecture and pre-training of Saprot. **c**, Saprot supports numerous protein prediction tasks, see more specific tasks in Supplementary Table 5. **d**, Performance comparison on 14 diverse prediction tasks. **e**, Saprot ranked *first* on ProteinGym leaderboard, outperforming 60 single and hybrid methods. Hybrid methods are not shown here (see Supplementary Fig. 1). **f**, Comparison of Saprot and ESM-2 in structural representation ability. The y-axis represents the cosine similarity score of the last layer representation for every pair of proteins, while the x-axis represents the TM-score of the protein pairs. Saprot shows higher Pearson correlation coefficient (ρ) between x and y axes, indicating more accurate protein structure encoding. **g**, Performance on protein sequence design. Although Saprot was not trained with a typical generation loss, its protein design performance rivals that of ProteinMPNN, with over 16 times faster inference speed. Saprot (AF2+CATH) indicates that it was pre-trained using AlphaFold2 predicted structures and fine-tuned with experimental structures from the CATH dataset. CATH is also used for training ProteinMPNN.

models through collective efforts. SaprotHub, as the first PLM Hub to join OPMC, marks the first step to bring the OPMC vision to fruition.

SaprotHub comprises three novel contributions: a new protein language model called Saprot (Fig. 1a,b), ColabSaprot (Fig. 2a,b,c) for fine-tuning (or re-training) and inference of Saprot using adapter learning techniques [17] on the Colab platform, and a centralized community model repository (Fig. 2d) for saving, sharing, searching, and collaborating on fine-tuned Saprot models. Through the innovative integration of

universal protein language models, Colab-based cloud computing, and adapter-based fine-tuning, SaproHub addresses several key challenges. These challenges include the difficulty of sharing and aggregating large-scale models, the risk of catastrophic forgetting during continual learning, and the need to protect proprietary data. Below, we describe these three modules in detail.

We developed Sapro, a cutting-edge, large-scale protein language model that forms the foundation for ColabSapro and SaproHub. A key feature of Sapro is the introduction of a novel structure-aware (SA) protein representation method (Fig. 1a), departing from classical amino acid (AA) sequences or explicit 3D structures. Specifically, we introduce a unique SA vocabulary (*SAV*) where $SAV = \{\mathcal{V}, \mathcal{F}\}$, combining 20 amino acids (\mathcal{V}) and 20 structural (3Di) alphabet tokens (\mathcal{F}), generated by encoding 3D structures with Foldseek [39]. SA tokens encompass all AA and 3Di combinations in *SAV*, allowing proteins to be represented as SA strings that capture both primary and tertiary structures (Fig. 1b). Despite its conciseness, applying the SA vocabulary successfully addresses the key challenges of scalability and overfitting in training on large-scale *AlphaFold2-generated* atomic structures (see [37, 11, 4]). The elegant use of SA token sequences to represent proteins has gained increasing attention in subsequent studies [9, 24, 33, 11, 19] and holds promise as a new paradigm for protein representation.

We adopt a bidirectional Transformer [47] for Sapro (Fig. 1b). During pre-training, Sapro uses SA sequences instead of AA sequences, *partially* masks random SA tokens, and predicts their complete forms (Methods and Fig. 1b). Unlike models built on existing PLMs, Sapro was trained from scratch with 650 million parameters on over 40 million protein SA sequences, filtered from 214 million AlphaFold database proteins [40] with a 50% identity threshold, over three months using 64 NVIDIA 80G A100 GPUs. This made it the largest structure dataset for PLM training at the time [37]. This huge computational demand mirrors that of ESM-2 (650M) and AlphaFold3 [1].

After pre-training, Sapro has become a foundational protein language model, excelling in various protein prediction tasks (Fig. 1c), including supervised training, zero-shot prediction, and protein sequence design. In supervised training, it handles both protein-level (e.g., fluorescence, thermostability, structure, and folding class prediction) and residue-level (e.g., binding site detection) tasks [18]. In zero-shot scenarios, Sapro predicts enzyme activity [12], viral mutation fitness [14], disease variant effects, and driver mutations [8] without fine-tuning (Methods). Though not designed with an explicit generative loss, Sapro (using masked language modeling loss) shows remarkable proficiency in protein sequence design, achieving a 16 \times speedup in inference compared to ProteinMPNN [3] (Fig. 1g, 2h, Supplementary Fig. 2).

Fig. 1d shows Sapro's exceptional performance across 14 diverse protein prediction tasks compared to the ESM-2 [25] and ProtBert [6] PLMs. When protein structures are available, Sapro consistently outperforms both models in tasks marked in purple, blue, and red. Without structures available, Sapro remains highly competitive (green tasks). Notably, Sapro with frozen representation substantially outperforms ESM-2 in three zero-shot mutational effect prediction tasks: Mega-scale [38] (0.574 vs. 0.478), ProteinGym [32] (0.457 vs. 0.414), and ClinVar [23] (0.909 vs. 0.862).

Sapro held the top spot on the public blind ProteinGym leaderboard (Fig. 1e, Supplementary Fig. 1) for several months after its release. It surpasses over 60 notable methods, including ESM-2 (650M, 3B, 15B), EVE [8], ESM-1v [54], MSA Transformer [59], ESM-IF [15], DeepSequence [35], and all hybrid models. Notably, the 35M parameter version of Sapro outperforms the 15B parameter version of ESM-2. Supplementary Tables 2-3 and Fig. 1, 3 provide more comparisons. Recent studies also show Sapro's impressive results in B-cell conformational epitope prediction [19], protein engineering [42] and drug-target interaction [28], etc. These results show that Sapro has an exceptional 'one-model-fits-all' capability (Fig. 1c, d and Supplementary Table 5), critical to achieve SaproHub's community collaboration vision.

We developed ColabSapro by integrating Sapro into Google Colab's infrastructure to support PLM training. ColabSapro enables seamless deployment and execution of various task-specific re-trained Sapro models, eliminating the need for environment setup and code debugging. It also allows researchers to initiate training sessions with just a few clicks (Fig. 2a). ColabSapro is designed to accommodate all tasks within the original Sapro framework, enabling direct prediction for tasks such as zero-shot mutation effect prediction [12, 8, 35] and protein sequence design [16, 3]. For mutation effect prediction, it implements single-site, multi-site, and whole protein sequence single-site saturation mutation. For sequence design, it can design sequences based on backbone structures or redesign specific protein regions, such as antibody variable regions, while keeping the scaffold fixed (Methods).

For supervised re-training tasks, users can fine-tune ColabSapro with their own data. Here, we introduce a parameter-efficient fine-tuning technique by integrating lightweight adapter networks into ColabSapro (Fig. 2b). During training, only adapter parameters are updated, achieving nearly the same accuracy as fine-tuning all Sapro parameters (Supplementary Fig. 4). This design not only enhances optimization efficiency [41, 63], but more importantly, as proposed here for the first time, it enables collaborative and centralized mechanisms for these biologist fine-tuned, task-specific PLMs (or Sapro) within the community, especially in the cloud environment (Fig. 2c, d, e, f and Methods). With adapters, researchers can also easily store and exchange their re-trained Sapro models in SaproHub by loading or uploading adapter networks instead of the full pre-trained model. Since adapter networks contain much fewer parameters (around 1% of the whole Sapro model), this method greatly reduces storage, communication, and management burdens

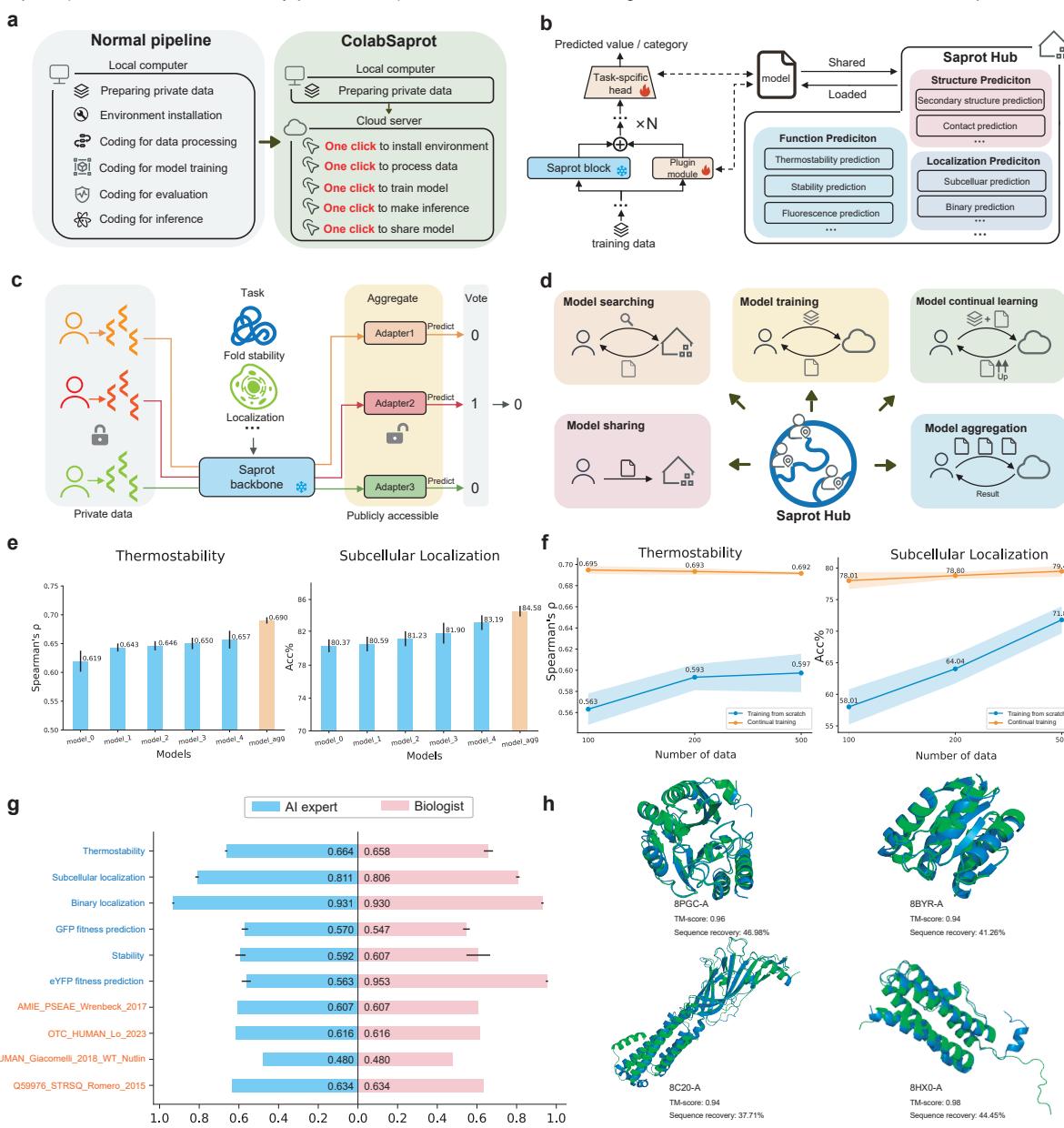


Fig. 2: Illustration of ColabSaprot and SaprotHub. a, Comparison of model training and inferring between normal pipeline and ColabSaprot. ColabSaprot streamlines the process, offering simplified procedures for model training and inference with just a few clicks. **b,** A lightweight plugin module (i.e., adapter) is inserted into Saprot for efficient training, sharing and co-construction. **c,** ColabSaprot predicts protein functions by aggregating existing models (or adapters) in SaprotHub without requiring private data. **d,** SaprotHub enables biologists to co-share, co-build, co-use and collaborate within the community. **e,** ColabSaprot’s community-wide model collaboration mechanism (see **c**) allows it to achieve higher performance (orange bars) by aggregating individually trained models (blue bars). Each individual model is trained with its own data, which may or may not overlap. **f,** By continually learning on models trained and shared by other biologists, ColabSaprot substantially outperforms those training-from-scratch models, especially when users lack sufficient training data (x-axis represents the number of training examples). **g,** User study on supervised fine-tuning and zero-shot mutation effect prediction tasks. **h,** User study on the inverse folding task. Experimental structures are shown in green, while predicted structures are in blue.

on SaprotHub, making accessibility, co-construction, co-utilization, and co-sharing possible for a broader scientific community (Fig. 2b, c, d).

Additionally, we developed many features for ColabSaprot. Our interface supports nine types of data inputs, including SA sequences, AA sequences, UniProt IDs, and PDB files, etc. We offer automatic conversion of input formats, division of training and testing sets, overfitting checks, automatic best checkpoint saving, evaluation, and visualization. To address GPU memory limitations, we implemented an adaptive batch size feature based on gradient accumulation [22]. We also provide advanced settings that allow researchers to extend code and customize training hyperparameters to suit their specific needs.

123 Finally, we developed Saprothub, a centralized model store for seamlessly sharing and co-building peer-
124 retrained Saproth models within the biology community (Fig. 2d). Through the ColabSaproth interface, we
125 implemented features like model storage, sharing, search, continual learning, and aggregation for better per-
126 formance. This brings two key benefits to Saprothub: (1) Biologists can share their trained models without
127 compromising sensitive data, overcoming data sharing challenges and enabling the sharing of research out-
128 comes (Fig. 2c); (2) Biologists can continually train models based on shared models from peers, especially
129 when data is limited, as fine-tuning on a superior initialization model typically improves predictive accuracy
130 (Fig. 2f). As more biologists join Saprothub, it will accumulate a large collection of models for various pro-
131 tein prediction tasks. Additionally, multiple Saproth models will be available for the same protein function.
132 New users can enhance predictive performance (Fig. 2e) by building upon existing models through adapter
133 aggregation (Fig. 2c and Methods) or advanced ensemble learning [21].

134 Saprothub has attracted dozens of OPMC developers, sharing over 30 task-specific models. Among
135 these, the Zheng team submitted a high-quality eYFP fluorescence fitness prediction model with a Spearman
136 correlation (ρ) of 0.94, nearing wet lab accuracy on double-site and triple-site mutation tasks (see Methods:
137 Saprothub page, model name: Model-EYFP_100K-650M). Trained on the largest proprietary fluorescent
138 protein mutation dataset of 100,000 experimentally-tested variants, this model's impressive accuracy is
139 crucial for advancing fluorescent protein design and applications in biotechnology, medical research and
140 evolutionary biology. Additionally, the Chang team recently utilized ColabSaproth to conduct zero-shot single-
141 point mutation predictions on a uracil-N-glycosylase variant, eTDG, and experimentally validated its T-to-G
142 editing efficiency. Among the top 20 predicted mutations, 17 exhibited higher editing efficiency than the
143 wild type, with L74E, H11K, and L74Q showing nearly a 2-fold increase (see Methods and Supplementary
144 Table 6).

145 We also conducted a user study by recruiting 12 biology researchers (without an ML background) and
146 compared their performance with that of an AI expert (Methods). The results demonstrated that, with
147 ColabSaproth and Saprothub, biology researchers can train and use state-of-the-art protein models with
148 performance comparable to that of an AI expert (Fig. 2g,h). In some scenarios, such as the eYFP fitness
149 prediction task, biologists using existing models in Saprothub achieved notably higher predictive accuracy
150 than AI experts, particularly when the training dataset was unavailable or limited.

151 Saprothub and ColabSaproth empower biology researchers to easily train and share advanced protein
152 ML models for various prediction tasks, even without expertise in ML or coding. This capability enables
153 the entire protein research community to contribute models, fostering collaboration and promoting the
154 dissemination of high-quality, peer-trained models. We have open-sourced the code for both Saproth and
155 ColabSaproth (Methods), allowing other PLMs to adopt a similar approach to build their own PLM Hub or
156 join Saprothub. Our OPMC members are continuously integrating more PLMs (such as ESM, ProtTrans,
157 etc.) into the OPMC framework, vastly expanding access to a diverse array of PLMs

158 This community-wide participation approach to protein modeling aligns with the OPMC vision. Our goal
159 here is to inspire and foster the cooperative construction of open protein prediction models through Sapro-
160 thub. We envision Saprothub as the catalyst that initializes OPMC, driving innovation and collaboration
161 in the field.

162 References

- 163 [1] Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ron-
164 neberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure prediction
165 of biomolecular interactions with alphafold 3. *Nature*, pages 1–3, 2024.
- 166 [2] Ethan C Alley, Grigory Khimulya, Surojit Biswas, Mohammed AlQuraishi, and George M Church.
167 Unified rational protein engineering with sequence-based deep representation learning. *Nature methods*,
168 16(12):1315–1322, 2019.
- 169 [3] Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F Milles,
170 Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep learning-based
171 protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.
- 172 [4] Diego del Alamo. The issues when training on the afdb structure data. *twitter*:
173 <https://x.com/DdelAlamo/status/1795353297580445851>, 2024.
- 174 [5] Jurgen Drews. Drug discovery: a historical perspective. *science*, 287(5460):1960–1964, 2000.
- 175 [6] Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom
176 Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, et al. Prottrans: Toward understanding the
177 language of life through self-supervised learning. *IEEE transactions on pattern analysis and machine
178 intelligence*, 44(10):7112–7127, 2021.

- 179 [7] Noelia Ferruz, Steffen Schmidt, and Birte Höcker. Protgpt2 is a deep unsupervised language model for
180 protein design. *Nature communications*, 13(1):4348, 2022.
- 181 [8] Jonathan Frazer, Pascal Notin, Mafalda Dias, Aidan Gomez, Joseph K Min, Kelly Brock, Yarin Gal,
182 and Debora S Marks. Disease variant prediction with deep generative models of evolutionary data.
183 *Nature*, 599(7883):91–95, 2021.
- 184 [9] Benoit Gaujac, Jérémie Donà, Liviu Copoiu, Timothy Atkinson, Thomas Pierrot, and Thomas D
185 Barrett. Learning the language of protein structure. *arXiv preprint arXiv:2405.15840*, 2024.
- 186 [10] Michael H Glickman and Aaron Ciechanover. The ubiquitin-proteasome proteolytic pathway:
187 destruction for the sake of construction. *Physiological reviews*, 2002.
- 188 [11] Thomas Hayes, Roshan Rao, Halil Akin, Nicholas James Sofroniew, Deniz Oktay, Zeming Lin, Robert
189 Verkuil, Vincent Quy Tran, Jonathan Deaton, Marius Wiggert, Rohil Badkundri, Irhum Shafkat, Jun
190 Gong, Alexander Derry, Raul Santiago Molina, Neil Thomas, Yousuf Khan, Chetan Mishra, Carolyn
191 Kim, Liam J. Bartie, Patrick D. Hsu, Tom Sercu, Salvatore Candido, and Alexander Rives. Simulating
192 500 million years of evolution with a language model. *bioRxiv*, 2024.
- 193 [12] Yan He, Xibin Zhou, Chong Chang, Ge Chen, Weikuan Liu, Geng Li, Xiaoqi Fan, Mingsun Sun, Chensi
194 Miao, Qianyue Huang, et al. Protein language models-assisted optimization of a uracil-n-glycosylase
195 variant enables programmable t-to-g and t-to-c base editing. *Molecular Cell*.
- 196 [13] Michael Heinzinger, Konstantin Weissenow, Joaquin Gomez Sanchez, Adrian Henkel, Martin Steinegger,
197 and Burkhard Rost. Prostt5: Bilingual language model for protein sequence and structure. *bioRxiv*,
198 pages 2023–07, 2023.
- 199 [14] Brian Hie, Ellen D Zhong, Bonnie Berger, and Bryan Bryson. Learning the language of viral evolution
200 and escape. *Science*, 371(6526):284–288, 2021.
- 201 [15] Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander
202 Rives. Learning inverse folding from millions of predicted structures. In *International conference on
203 machine learning*, pages 8946–8970. PMLR, 2022.
- 204 [16] Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander
205 Rives. Learning inverse folding from millions of predicted structures. In *International conference on
206 machine learning*, pages 8946–8970. PMLR, 2022.
- 207 [17] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
208 Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*,
209 2021.
- 210 [18] Mingyang Hu, Fajie Yuan, Kevin Yang, Fusong Ju, Jin Su, Hui Wang, Fei Yang, and Qiuyang Ding.
211 Exploring evolution-aware &-free protein language models as protein function predictors. *Advances in
212 Neural Information Processing Systems*, 35:38873–38884, 2022.
- 213 [19] Nikita V Ivanisenko, Tatiana I Shashkova, Andrey Shevtsov, Maria Sindeeva, Dmitriy Umerenkov, and
214 Olga Kardymon. Sema 2.0: web-platform for b-cell conformational epitopes prediction using artificial
215 intelligence. *Nucleic Acids Research*, page gkae386, 2024.
- 216 [20] François Jacob and Jacques Monod. Genetic regulatory mechanisms in the synthesis of proteins. *Journal
217 of molecular biology*, 3(3):318–356, 1961.
- 218 [21] Dongfu Jiang, Xiang Ren, and Bill Yuchen Lin. Llm-blender: Ensembling large language models with
219 pairwise ranking and generative fusion. *arXiv preprint arXiv:2306.02561*, 2023.
- 220 [22] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger,
221 Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate
222 protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- 223 [23] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth R Brown, Chen Chao, Shanmuga Chitipiralla,
224 Baoshan Gu, Jennifer Hart, Douglas Hoffman, Wonhee Jang, et al. Clinvar: improving access to variant
225 interpretations and supporting evidence. *Nucleic acids research*, 46(D1):D1062–D1067, 2018.

- 226 [24] Mingchen Li, Yang Tan, Xinzhu Ma, Bozitao Zhong, Huiqun Yu, Ziyi Zhou, Wanli Ouyang, Bingxin
227 Zhou, Liang Hong, and Pan Tan. Prosst: Protein language modeling with quantized structure and
228 disentangled attention. *bioRxiv*, pages 2024–04, 2024.
- 229 [25] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert
230 Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure
231 with a language model. *Science*, 379(6637):1123–1130, 2023.
- 232 [26] Ali Madani, Bryan McCann, Nikhil Naik, Nitish Shirish Keskar, Namrata Anand, Raphael R Eguchi,
233 Po-Ssu Huang, and Richard Socher. Progen: Language modeling for protein generation. *arXiv preprint arXiv:2004.03497*, 2020.
- 235 [27] Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alex Rives. Language models enable zero-shot prediction of the effects of mutations on protein function. *Advances in neural information processing systems*, 34:29287–29303, 2021.
- 238 [28] Zhaohan Meng, Zaiqiao Meng, and Iadh Ounis. Fusiondti: Fine-grained binding discovery with token-level fusion for drug-target interaction. *arXiv preprint arXiv:2406.01651*, 2024.
- 240 [29] Milot Mirdita, Konstantin Schütze, Yoshitaka Moriwaki, Lim Heo, Sergey Ovchinnikov, and Martin
241 Steinegger. Colabfold: making protein folding accessible to all. *Nature methods*, 19(6):679–682, 2022.
- 242 [30] Erik Nijkamp, Jeffrey A Ruffolo, Eli N Weinstein, Nikhil Naik, and Ali Madani. Progen2: exploring
243 the boundaries of protein language models. *Cell systems*, 14(11):968–978, 2023.
- 244 [31] Pascal Notin, Mafalda Dias, Jonathan Frazer, Javier Marchena-Hurtado, Aidan N Gomez, Debora
245 Marks, and Yarin Gal. Tranception: protein fitness prediction with autoregressive transformers and
246 inference-time retrieval. In *International Conference on Machine Learning*, pages 16990–17017. PMLR,
247 2022.
- 248 [32] Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood Van Niekerk, Steffanie Paul, Han Spinner, Nathan
249 Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, et al. Proteingym: large-scale benchmarks for
250 protein fitness prediction and design. *Advances in Neural Information Processing Systems*, 36, 2024.
- 251 [33] Mahdi Pourmirzaei, Farzaneh Esmaili, Mohammadreza Pourmirzaei, Duolin Wang, and Dong Xu.
252 Prot2token: A multi-task framework for protein language processing using autoregressive language
253 modeling. *bioRxiv*, pages 2024–05, 2024.
- 254 [34] Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and
255 Alexander Rives. Msa transformer. In *International Conference on Machine Learning*, pages 8844–8856.
256 PMLR, 2021.
- 257 [35] Adam J Riesselman, John B Ingraham, and Debora S Marks. Deep generative models of genetic
258 variation capture the effects of mutations. *Nature methods*, 15(10):816–822, 2018.
- 259 [36] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle
260 Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from scaling unsupervised
261 learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*,
262 118(15):e2016239118, 2021.
- 263 [37] Jin Su, Chenchen Han, Yuyang Zhou, Junjie Shan, Xibin Zhou, and Fajie Yuan. Saprot: Protein language
264 modeling with structure-aware vocabulary. In *The Twelfth International Conference on Learning Representations*, 2023.
- 266 [38] Kotaro Tsuboyama, Justas Dauparas, Jonathan Chen, Elodie Laine, Yasser Mohseni Behbahani,
267 Jonathan J Weinstein, Niall M Mangan, Sergey Ovchinnikov, and Gabriel J Rocklin. Mega-scale
268 experimental analysis of protein folding stability in biology and design. *Nature*, 620(7973):434–444,
269 2023.
- 270 [39] Michel van Kempen, Stephanie S Kim, Charlotte Tumescheit, Milot Mirdita, Cameron LM Gilchrist,
271 Johannes Söding, and Martin Steinegger. Foldseek: fast and accurate protein structure search. *Biorxiv*,
272 pages 2022–02, 2022.
- 273 [40] Mihaly Varadi, Stephen Anyango, Mandar Deshpande, Sreenath Nair, Cindy Natassia, Galabina Yordanova,
274 David Yuan, Oana Stroe, Gemma Wood, Agata Laydon, et al. AlphaFold protein structure

275 database: massively expanding the structural coverage of protein-sequence space with high-accuracy
276 models. *Nucleic acids research*, 50(D1):D439–D444, 2022.

277 [41] Shuai Zeng, Duolin Wang, and Dong Xu. Peft-sp: Parameter-efficient fine-tuning on large protein
278 language models improves signal peptide prediction. *bioRxiv*, pages 2023–11, 2023.

279 [42] Ziyi Zhou, Liang Zhang, Yuanxi Yu, Banghao Wu, Mingchen Li, Liang Hong, and Pan Tan. Enhancing
280 efficiency of protein language models with minimal wet-lab data through few-shot learning. *Nature
281 Communications*, 15(1):5566, 2024.

282 OPMC (senior) member list

283 Sergey Ovchinnikov, MIT
284 Martin Steinegger, Seoul National University
285 Kevin K. Yang, Microsoft
286 Michael Heinzinger, Technische Universität München
287 Pascal Notin, Harvard University
288 Pranam Chatterjee, Duke University
289 Jia Zheng, Westlake University
290 Stan Z. Li, Westlake University
291 Xing Chang, Westlake University
292 Huaizong Shen, Westlake University
293 Noelia Ferruz, The Centre for Genomic Regulation (CRG)
294 Rohit Singh, Duke University
295 Debora Marks, Harvard University
296 Anping Zeng, Westlake University
297 Jijie Chai, Westlake University
298 Anthony Gitter, University of Wisconsin-Madison
299 Anum Glasgow, Columbia University
300 Milot Mirdita, Seoul National University
301 Philip M. Kim, University of Toronto
302 Christopher Snow, Colorado State University
303 Vasilis Ntranos, University of California
304 Philip A. Romero, Duke University
305 Jianyi Yang, Shandong University
306 Caixia Gao, Chinese Academy of Sciences
307

308 OPMC (regular) member list

309 Jiawei Zhang, Westlake University
310 Yuyang Tao, ShanghaiTech University
311 Fengyuan Dai, Westlake University
312 Xuting Zhang, Westlake University
313 Yan He, Westlake University
314
315 New authors will be added until the final accepted version.

316 Methods

317 Constructing SA Protein Sequence

318 The Structure-Aware (SA) vocabulary encompasses both residue and structure information, as illustrated in
319 Fig. 1a. Given a protein P , its primary sequence can be denoted as (s_1, s_2, \dots, s_n) , where $s_i \in \mathcal{V}$ represents the
320 residue at the i_{th} site, and \mathcal{V} represents residue alphabet. Drawing inspiration from the vector quantization
321 learning technique [64], we encode protein tertiary structures into discrete residue-like structural tokens.
322 Here, we use Foldseek [39], a fast and accurate protein structure aligner. Through Foldseek, we have a
323 structure alphabet \mathcal{F} , wherein P is expressed as the sequence (f_1, f_2, \dots, f_n) , with $f_j \in \mathcal{F}$ representing
324 the 3Di token for the j_{th} residue site. To maintain simplicity, we adhere to the default configuration of
325 Foldseek, which sets the size m of \mathcal{F} to 20. We then combine the residue and structure tokens per residue
326 site, generating a new structure-aware protein sequence $P = (s_1f_1, s_2f_2, \dots, s_nf_n)$, where $s_i f_i \in \mathcal{V} \times \mathcal{F}$ is the
327 so-called SA token naturally fusing both residue and geometric conformation information. The SA-token
328 protein sequence can then be fed into a standard Transformer encoder as basic input. It's important to note
329 that we also introduce a mask signal “#” to both residue and structure alphabet, which results in “ $s_i\#$ ”
330 and “ $\#f_i$ ” that indicate only residue or structure information is available. The size of the SA vocabulary is
331 $21 \times 21 = 441$.

332 Saprot architecture and pre-training

333 Saprot employs the same network architecture and parameter size as the 650M version of ESM-2 [25], which
334 is inspired by the BERT [47] model in the Natural Language Processing (NLP) field. The main distinction lies
335 in the expanded embedding layer, which encompasses 441 structure-aware tokens instead of the original 20
336 residue tokens. This nearly identical architecture enables straightforward comparison with the ESM model.

337 Saprot is pre-trained with the Masked Language Modeling (MLM) objective, similar to ESM-2 and
338 BERT. Formally, for a protein sequence $P = (s_1f_1, s_2f_2, \dots, s_nf_n)$, the input and output can be represented
339 as: $input : (s_1f_1, \dots, \#f_i, \dots, s_nf_n) \rightarrow output : s_i f_i$ (Fig. 1b). Given that the 3Di token may not always be
340 accurate for certain regions in predicted structures by AlphaFold2, f_i in $\#f_i$ is made visible during training
341 to reduce the emphasis the model places on its predictions.

342 Structures predicted by AlphaFold2 (AF2) are accompanied by confidence scores, referred to as pLDDT,
343 which provide an assessment of the precision of atom coordinates. Therefore, we have implemented special
344 handling for regions with low pLDDT scores. During the pre-training process, when regions with pLDDT
345 scores lower than 70 are chosen for MLM prediction, we predict the “ $s_i\#$ ” token. Here, the input SA
346 sequence is masked with the “##” token. This approach forces the model to focus on predicting the types
347 of residues in these regions based on their context. In cases where regions with lower pLDDT scores are
348 not selected for MLM prediction, the input still use the “ $s_i\#$ ” token. This ensures that the model solely
349 relies on residue context (instead of the structural context) in these regions to assist in predicting other
350 tokens. For downstream task prediction phases, we maintain consistency with the training data by handling
351 regions with lower pLDDT scores (still with 70 as threshold). Specifically, tokens within these regions are
352 represented as “ $s_i\#$ ”, with only residue tokens being visible.

353 Processing pre-training data

354 We followed the procedures outlined in ESM-2 [25] to generate filtered protein data. Subsequently, we
355 acquired all AF2 structures via the AlphaFoldDB website (<https://alphafold.ebi.ac.uk/>), utilizing the
356 UniProt IDs of protein sequences. Proteins without structures in AlphaFoldDB were omitted for pre-training.
357 This process yielded a collection of approximately 40 million structures. Employing Foldseek, we encoded
358 these structures into 3Di tokens. Subsequently, we formulated structure-aware sequences by combining
359 residue and 3Di tokens, treating them as a single SA token at each position.

360 Pre-training hyper-parameters

361 Following ESM-2 and BERT, during training, 15% of the SA tokens in each batch are masked. We replace
362 the SA token $s_i f_i$ with the $\#f_i$ token 80% of the time, while 10% of the tokens are replaced with randomly
363 selected tokens, and the other 10% tokens remain unchanged. For the optimization of Saprot, we adopt simi-
364 lar hyper-parameters to those employed in the ESM-2 training phase. Specifically, we employ the AdamW
365 optimizer [53], setting $\beta_1 = 0.9$, $\beta_2 = 0.98$ and we utilize L_2 weight decay of 0.01. We gradually increase the
366 learning rate from 0 to 4e-4 over the first 2000 steps and linearly lower it to 5e-4 from 150K steps to 1.5M
367 steps. The overall training phase lasts approximately 3M steps. Like the ESM model, we also truncate them
368 to a maximum of 1024 tokens, and our batch size consists of 512 sequences. Additionally, we also employ
369 mixed precision training to train Saprot.

370 Baseline models

371 We compared Saprot to several prominent protein language models (see Supplementary Table 2,3). For
372 supervised learning tasks, we compared Saprot with ESM-2 (the 650M version) [25], ProtBert (the BFD
373 version) [6], MIF-ST [67], and GearNet [69]. The first two models utilize residue sequences as input, while
374 the latter two models incorporate both residue and 3D structures as input. ESM-2 (650M) stands out as
375 the primary baseline for comparison, given its similar model architecture, size, and training approach when
376 compared to Saprot. ESM-2 also offers a 15 billion (B) parameter version, which can be challenging to fine-
377 tune even on GPUs with 80G memory. Therefore, we only conducted comparisons with ESM-2 (15B) for
378 zero-shot mutational effect prediction tasks, which can be achieved without the need for fine-tuning.

379 For the zero-shot mutational effect prediction task, we compare with the state-of-the-art ESM-2, Prot-
380 Bert, ESM-1v [54], Tranception L (without MSA retrieval) [55], MSA Transformer [59], and EVE [8] models.
381 For the protein inverse folding task, we compare with ProteinMPNN [46] as baseline (see Fig. 1).

382 Zero-shot formula

383 Previous sequence-based protein language models like the ESM models predict mutational effects using the
384 log odds ratio at the mutated position. The calculation can be formalized as follows:

$$\sum_{t \in T} [\log P(x_t = s_t^{mt} | x_{\setminus T}) - \log P(x_t = s_t^{wt} | x_{\setminus T})] \quad (1)$$

385 Here T represents all mutations and $s_t \in \mathcal{V}$ is the residue type for mutant and wild-type sequence. We
386 slightly modify the formula above to adapt to the structure-aware vocabulary, where the probability assigned
387 to each residue corresponds to the summation of tokens encompassing that specific residue type, as shown
388 below:

$$\sum_{t \in T} [\log \sum_{f \in \mathcal{F}} P(x_t = s_t^{mt} f | x_{\setminus T}) - \log \sum_{f \in \mathcal{F}} P(x_t = s_t^{wt} f | x_{\setminus T})] \quad (2)$$

389 Here $f \in \mathcal{F}$ is the 3Di token generated by Foldseek and $s_t f \in \mathcal{V} \times \mathcal{F}$ is the SA token in our new vocabulary.

390 Zero-shot benchmarks

391 **ProteinGym** [32] is an extensive set of Deep Mutational Scanning (DMS) assays, enabling thorough com-
392 parison among zero-shot predictors. We assess all baselines on the substitution branch of the ProteinGym,
393 utilizing its provided protein structures and following the standard evaluation pipeline from the original
394 paper. For evaluation, We adopt Spearman's rank correlation as our metric.

395 **ClinVar** [23] serves as a freely accessible and publicly available repository containing information about
396 human genetic variants and interpretations of their significance to disease. In our analysis, we harness the
397 data sourced from EVE [8], and filter out proteins with length greater than 1024. To enhance the reliability of
398 our data, we opt to consider proteins with labels 1 "Gold Stars" or higher, which indicate higher credibility.
399 Following the methodology employed in EVE, we evaluate models' performance using the AUC metric.

400 **Mega-scale** [38] uses cDNA display proteolysis to measure the thermodynamic folding stability of
401 protein domains. Its released dataset contains all single mutations and selected double mutants of natural
402 and *de novo* designed proteins. We employ Spearman's rank correlation as our metric.

403 For each mutation dataset, we provide all variants with the wild-type structure, as AF2 does not reli-
404 ably distinguish the structural changes induced by single mutations. The ClinVar dataset only provides
405 UniProt IDs, so we manually downloaded all AF2 structures and eliminated proteins without structures in
406 AlphaFoldDB. Both ProteinGym and Mega-scale datasets provide protein structures, either predicted from
407 AF2 or derived from *de novo* design.

408 Supervised fine-tuning benchmarks

409 **Fine-tuning Saprot with AF2 structure.** For benchmarks lacking PDB structures, we retrieve all
410 AF2-predicted structures using UniProt IDs. These benchmarks encompass the Thermostability task from
411 FLIP [45] and the Localization Prediction task from DeepLoc [43]. DeepLoc consists of two branches for
412 subcellular localization prediction: one involving 10 location categories and the other involving binary
413 localization prediction with 2 location categories. We assess the models' performance on both branches.

414 **Fine-tuning Saprot with PDB structure.** These tasks provide experimentally determined structures
415 as training data. We evaluate Metal Ion Binding task [50] and a series of tasks from ProteinShake [52],
416 including Structure Class Prediction, Structural Similarity Prediction and Binding Site Detection.

417 **Fine-tuning Saprot without structure.** Saprot is specifically developed to incorporate protein struc-
418 tures as valuable input. However, it can still work even in scenarios where structures are not provided. In
419 such cases, the 3Di token is masked within the SA sequence. We evaluate the Fluorescence prediction and
420 Stability prediction tasks from the TAPE [58] benchmark, the AAV dataset from the FLIP [45] benchmark

421 and β -lactamase landscape prediction from the PEER [66] benchmark. For a given wild type protein, its
422 variants do not offer additional structures.

423 Dataset splits

424 In previous literature, a common practice was to partition datasets based on sequence identity. However, a
425 recent benchmark study, ProteinShake [52], mentioned that protein structures exhibit greater conservation
426 than protein sequences. They argue that dissimilar sequences in training and test sets may share similar
427 structures, leading to data leakage issues. Consequently, we adopt a more strict data splitting approach
428 proposed by ProteinShake and split the data based on structure similarity. For all datasets with structures,
429 we employ a 70% structure similarity threshold for splitting. For those without provided structures, we still
430 follow the original data splits.

431 Supervised fine-tuning hyper-parameters

432 We employ the AdamW [53] optimizer during fine-tuning, setting $\beta_1 = 0.9$, $\beta_2 = 0.98$, along with L_2 weight
433 decay of 0.01. We use a batch size of 64 for all datasets. We empirically found that the optimal learning rate
434 for most baselines are in the range of 1e-5 to 5e-5. For **Training with AF2 structure** and **Training with**
435 **PDB structure**, the optimal learning rate was set to 5e-5, whereas for **Training without structure**, it
436 was set to 1e-5. We fine-tuned all model parameters until convergence and selected the best checkpoints
437 based on their performance on the validation set.

438 Protein inverse folding

439 Given the 3D coordinates of a protein backbone, protein inverse folding aims to predict residue sequences that
440 fold into this shape. This exhibits promising applications in protein *de novo* design. While not particularly
441 targeted on this task, Saprot can predict reliable residue sequences from the backbone due to the large-scale
442 pre-training tailored for predicting masked residues based on sequence and structure contexts. To utilize
443 Saprot for inverse folding, we first encode protein backbone into 3Di tokens (f_1, f_2, \dots, f_n). Then we add
444 masks to all residue parts of SA tokens, forming a SA sequence (# $f_1, #f_2, \dots, #f_n$), which serves as input
445 for Saprot to predict residue distributions at all positions. In contrast to ProteinMPNN [46] that generates
446 residues in the auto-regressive manner, i.e. generating next token conditioned on all previous outputs, Saprot
447 is able to simultaneously predict all residues with only one forward propagation. As shown in Fig. 1e, in
448 CATH [56] test set, Saprot predicts 16 times faster than ProteinMPNN, while achieves comparable prediction
449 accuracy, which demonstrates its huge potential in protein *de novo* design.

450 Embedding visualization

451 We employ t-SNE [65] to visualize the protein representations generated by Saprot and ESM-2 in three
452 datasets (Supplementary Fig. 3). We use the non-redundant version ($PIDE < 40\%$) of the SCOPe [44]
453 database, visualizing all alpha and beta proteins. We additionally visualize proteins in Subcellular Localization
454 and Binary Localization datasets from DeepLoc [43]. To generate protein representations, we adopt
455 the average pooling on the embeddings of the last layer for both Saprot and ESM-2.

456 Adapter learning

457 The adapter learning technique originated in the field of NLP [49, 48, 68] and has recently garnered attention
458 in protein-related research. However, these recent preliminary preprints primarily focused on its ability in
459 model accuracy or efficient GPU memory consumption [63, 62, 61, 60]. Through the application of adapters,
460 researchers have observed their ability to alleviate overfitting issues or integrate domain-specific knowledge
461 into protein language models. In this paper, we innovatively integrate it with Google Colaboratory and
462 endow it with deeper capabilities, namely, to promote model sharing and co-construction within the biologist
463 community in the online SaprotHub. This represents a novel perspective on the application of adapter
464 learning within protein biology. In the broader AI community, various types of adapters exist, including
465 Houldby [49], Pfeiffer [57], Compacter [51], and LoRA [17]. We choose LoRA for its capacity to deliver
466 comparable results with fewer parameters. By integrating learnable low-rank matrices into each Saprot
467 block while freezing the backbone, LoRA enables parameter-efficient fine-tuning and model sharing for these
468 downstream tasks.

469 Community-wide collaboration

470 As the SaprotHub community expands, an increasing number of Saprot models will become available for the
471 same protein function. This paves the way for collaborative co-construction among biologists. To illustrate
472 this advantage (Fig. 2e), we employ intuitive approach. In regression tasks, the final prediction is obtained
473 by calculating the mean value of outputs from all models. For classification tasks, we aggregate predicted

474 categories from all models and select the category with the highest prediction as the final outcome. Note
475 that all models shown in Fig. 2e were run five times using different random seeds. The lightweight adapter
476 technique plays a vital role in model aggregation, addressing the impracticality of loading and sharing
477 multiple large pre-trained models.

478 Furthermore, more advanced functionalities can be developed. For instance, by injecting multiple
479 adapters into Saprot via SaprotHub, continual learning can be conducted to further enhance performance.
480 Adapters eliminate the need to modify the weights of the original Saprot model, effectively mitigating the
481 issue of catastrophic forgetting. This study introduces a pioneering concept of integrating community wis-
482 dom by utilizing shared adapters, thereby ensuring data privacy. This design enables the realization of the
483 OPMC vision.

484 ColabSaprot notebooks

485 ColabSaprot comprises three key components. The first component focuses on training and inference,
486 enabling researchers to swiftly set up the runtime environment, preprocess training data, and fine-tune
487 Saprot with ease. We offer advanced hyper-parameter settings, allowing researchers to customize their train-
488 ing strategy based on specific requirements, such as batch size, learning rate, and training steps. After
489 completing the training process, researchers can share the model weights on the SaprotHub community
490 store, benefiting other biologists.

491 Furthermore, researchers without their own training data can directly access a range of fine-tuned Saprot
492 models from SaprotHub and perform inference for desired protein functions. These existing models also
493 support continual training with additional data, leading to improved prediction accuracy.

494 The second component focuses on the zero-shot mutational effect prediction task, enabling users to
495 leverage Saprot for predicting fitness changes resulting from residue mutations. We provide multiple options
496 for users to predict single mutations, multi-site mutations, or saturation mutations on given proteins.

497 The final component of Saprot offers support for the protein inverse folding task based on the backbone
498 structure. Saprot achieves a 16 times faster inference speed than ProteinMPNN while maintaining compet-
499 itive prediction performance. This capability empowers users to conduct precise protein designs on a large
500 scale.

501 User study

502 We recruited 12 participants to use ColabSaprot for 8 protein prediction tasks, including supervised fine-
503 tuning, zero-shot mutation effect prediction, and protein inverse folding. All these participants had an
504 education background in biology but no experience with machine learning. They were required to complete
505 their assigned tasks within 3 days. In comparison, the AI expert is a third-year PhD student in machine
506 learning with over 2 years of research experience in AI for protein studies. Our AI expert directly used the
507 Saprot code from our GitHub repository to perform these tasks. The AI expert conducted these experiments
508 five times using different random seeds with optimal hyper-parameters.

509 For the supervised fine-tuning, we utilized five publicly available datasets dataset, including those for
510 thermostability prediction, (subcellular and binary) localization prediction, GFP fluorescence prediction,
511 and stability prediction. Additionally, we introduced a proprietary eYFP fitness prediction dataset collected
512 by Zheng's lab, which has 3087 and 3088 samples for validation and testing. To reduce training time and
513 computing power consumption, we randomly selected 1,000 samples from the training set of each dataset
514 while keeping the validation and test sets unchanged.

515 Among the 12 participants, 10 conducted the supervised fine-tuning tasks. Each of these participants was
516 assigned three tasks from a set of six. In other words, each training task was performed by 5 participants.
517 We evaluated their average accuracy.

518 These participants were required to train the models and provide predictions on the corresponding test
519 sets for subsequent evaluation. It is important to note that for the eYFP task, biology participants were
520 required to perform continual learning on a model that had previously been trained and shared on the
521 SaprotHub platform (see <https://huggingface.co/SaProtHub/Model-EYFP-650M>), while the AI expert was
522 required to use the base Saprot model after pre-training. The purpose here is to demonstrate the advantage
523 of SaprotHub for researchers who lack training data. For all other tasks, biology participants and the AI
524 expert trained their models using the same base Saprot model (i.e., https://huggingface.co/westlake-repl/SaProt_650M_AF2) for fair evaluation.

526 For the zero-shot mutational effect prediction, we randomly selected four mutation datasets from Prote-
527 inGym benchmark [32] for evaluation. Three of these datasets focus on the impact of mutations on enzyme
528 activity, while the fourth addresses drug resistance. We assigned one participant to conduct predictions on
529 these datasets using ColabSaprot in a zero-shot manner.

530 For the protein inverse folding task, we assigned one participant to use ColabSaprot to generate protein
531 sequences based on given structures. Subsequently, the participant employed ESMFold (an interface provided
532 by ColabSaprot) to predict the structures of the generated sequences. To assess Saprot's ability on new
533 structures, we selected these recently released structures (see Fig. 2h).

534 Please found the user study materials at *Data availability* section below.

535 Wet-lab experimental validation on eTDG

536 The Chang team utilized ColabSaprot to perform zero-shot mutation effect prediction on a variant of
537 uracil-N-glycosylase called eTDG. They then conducted experimental validation following the procedure
538 outlined in [12] to verify the top 20 predicted mutations by ColabSaprot. Supplementary Table 6 provides
539 the predicted scores and experimentally validated editing efficiencies. Notably, the original paper by [12]
540 reported only 8 of the variants predicted by ColabSaprot.

541 Data availability

542 The pre-training dataset for training Saprot is available at <https://huggingface.co/datasets/westlake-repl/>
543 AF2_UniRef50. Downstream task datasets are all stored at <https://huggingface.co/SaProtHub>. Data for user
544 study analysis is available at [https://drive.google.com/file/d/1LdGRnwt2lttszNBAJ0F967A8rguPq8b/](https://drive.google.com/file/d/1LdGRnwt2lttszNBAJ0F967A8rguPq8b/view?usp=sharing)
545 view?usp=sharing.

546 Code and service availability

547 Saprot is an open-sourced model with MIT license. The code is available at <https://github.com/westlake-repl/Saprot>. The code implementation of ColabSaprot notebook is available at <https://github.com/westlake-repl/SaProtHub>. ColabSaprot service is available at <https://colab.research.google.com/github/westlake-repl/SaprotHub/blob/main/colab/SaprotHub.ipynb>. Fine-tuned models can be obtained through
548 SaprotHub <https://huggingface.co/SaProtHub>.

549 References

- 550 [43] José Juan Almagro Armenteros, Casper Kaae Sønderby, Søren Kaae Sønderby, Henrik Nielsen, and Ole
551 Winther. Deeploc: prediction of protein subcellular localization using deep learning. *Bioinformatics*,
552 33(21):3387–3395, 2017.
- 553 [44] John-Marc Chandonia, Naomi K Fox, and Steven E Brenner. SCOPe: classification of large macro-
554 molecular structures in the structural classification of proteins—extended database. *Nucleic Acids
555 Research*, 47(D1):D475–D481, 11 2018.
- 556 [45] Christian Dallago, Jody Mou, Kadina E Johnston, Bruce J Wittmann, Nicholas Bhattacharya, Samuel
557 Goldman, Ali Madani, and Kevin K Yang. Flip: Benchmark tasks in fitness landscape inference for
558 proteins. *bioRxiv*, pages 2021–11, 2021.
- 559 [46] Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J Ragotte, Lukas F Milles,
560 Basile IM Wicky, Alexis Courbet, Rob J de Haas, Neville Bethel, et al. Robust deep learning-based
561 protein sequence design using proteinmpnn. *Science*, 378(6615):49–56, 2022.
- 562 [47] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
563 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 564 [48] Junchen Fu, Fajie Yuan, Yu Song, Zheng Yuan, Mingyue Cheng, Shenghui Cheng, Jiaqi Zhang, Jie
565 Wang, and Yunzhu Pan. Exploring adapter-based transfer learning for recommender systems: Empirical
566 studies and practical insights. In *Proceedings of the 17th ACM International Conference on Web Search
567 and Data Mining*, pages 208–217, 2024.
- 568 [49] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea
569 Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In
570 *International conference on machine learning*, pages 2790–2799. PMLR, 2019.
- 571 [50] Mingyang Hu, Fajie Yuan, Kevin Yang, Fusong Ju, Jin Su, Hui Wang, Fei Yang, and Qiuyang Ding.
572 Exploring evolution-aware & -free protein language models as protein function predictors. In
573 S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural
574 Information Processing Systems*, volume 35, pages 38873–38884. Curran Associates, Inc., 2022.
- 575 [51] Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank
576 hypercomplex adapter layers. *Advances in Neural Information Processing Systems*, 34:1022–1035, 2021.

- 580 [52] Tim Kucera, Carlos Oliver, Dexiong Chen, and Karsten Borgwardt. Proteinshake: Building datasets
581 and benchmarks for deep learning on protein structures. *Advances in Neural Information Processing
582 Systems*, 36, 2024.
- 583 [53] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101,
584 2017.
- 585 [54] Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alex Rives. Language mod-
586 els enable zero-shot prediction of the effects of mutations on protein function. *Advances in neural
587 information processing systems*, 34:29287–29303, 2021.
- 588 [55] Pascal Notin, Mafalda Dias, Jonathan Frazer, Javier Marchena Hurtado, Aidan N Gomez, Debora
589 Marks, and Yarin Gal. Tranception: protein fitness prediction with autoregressive transformers and
590 inference-time retrieval. In *International Conference on Machine Learning*, pages 16990–17017. PMLR,
591 2022.
- 592 [56] Christine A Orengo, Alex D Michie, Susan Jones, David T Jones, Mark B Swindells, and Janet M
593 Thornton. Cath—a hierarchic classification of protein domain structures. *Structure*, 5(8):1093–1109,
594 1997.
- 595 [57] Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder,
596 Kyunghyun Cho, and Iryna Gurevych. Adapterhub: A framework for adapting transformers. *arXiv
597 preprint arXiv:2007.07779*, 2020.
- 598 [58] Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel,
599 and Yun Song. Evaluating protein transfer learning with tape. *Advances in neural information
600 processing systems*, 32, 2019.
- 601 [59] Roshan M Rao, Jason Liu, Robert Verkuil, Joshua Meier, John Canny, Pieter Abbeel, Tom Sercu, and
602 Alexander Rives. Msa transformer. In *International Conference on Machine Learning*, pages 8844–8856.
603 PMLR, 2021.
- 604 [60] Robert Schmirler, Michael Heinzinger, and Burkhard Rost. Fine-tuning protein language models boosts
605 predictions across diverse tasks. *bioRxiv*, pages 2023–12, 2023.
- 606 [61] Amelie Schreiber. Esmbind and qbind: Lora, qlora, and esm-2 for predicting binding sites and post
607 translational modification. *bioRxiv*, pages 2023–11, 2023.
- 608 [62] Samuel Sledzieski, Meghana Kshirsagar, Minkyung Baek, Bonnie Berger, Rahul Dodhia, and Juan Lav-
609 ista Ferres. Democratizing protein language models with parameter-efficient fine-tuning. *bioRxiv*,
610 2023.
- 611 [63] Samuel Sledzieski, Meghana Kshirsagar, Bonnie Berger, Rahul Dodhia, and Juan Lavista Ferres.
612 Parameter-efficient fine-tuning of protein language models improves prediction of protein-protein
613 interactions.
- 614 [64] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural
615 information processing systems*, 30, 2017.
- 616 [65] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine
617 Learning Research*, 9(86):2579–2605, 2008.
- 618 [66] Minghao Xu, Zuobai Zhang, Jiarui Lu, Zhaocheng Zhu, Yangtian Zhang, Ma Chang, Runcheng Liu,
619 and Jian Tang. Peer: a comprehensive and multi-task benchmark for protein sequence understanding.
620 *Advances in Neural Information Processing Systems*, 35:35156–35173, 2022.
- 621 [67] Kevin K Yang, Niccolò Zanichelli, and Hugh Yeh. Masked inverse folding with sequence transfer for
622 protein representation learning. *Protein Engineering, Design and Selection*, 36:gzad015, 2023.
- 623 [68] Fajie Yuan, Xiangnan He, Alexandros Karatzoglou, and Liguang Zhang. Parameter-efficient transfer
624 from sequential behaviors for user modeling and recommendation. In *Proceedings of the 43rd Interna-
625 tional ACM SIGIR conference on research and development in Information Retrieval*, pages 1469–1478,
626 2020.

- 627 [69] Zuobai Zhang, Minghao Xu, Arian Jamasb, Vijil Chenthamarakshan, Aurelie Lozano, Payel Das,
628 and Jian Tang. Protein representation learning by geometric structure pretraining. *arXiv preprint*
629 *arXiv:2203.06125*, 2022.

630 **Acknowledgements**

631 We thank Nan Li and Westlake University HPC Center for computing resources, and Jing Huang, Jianyang
632 Zeng, Dangshen Li, Longxing Cao, Dapeng Li for discussions and paper advice. We thank John M. Jumper
633 for providing insights about AF2 Evoformer and for his advice on the SA token design via email. This work
634 is supported by the National Key Research and Development Program of China (No. 2022ZD0115100),
635 the National Natural Science Foundation of China (No. U21A20427), the Westlake Center of Synthetic
636 Biology and Integrated Bioengineering (WE-SynBio), and the Research Center for Industries of the Future
637 (No. WU2022C030).

638 **Author contributions**

639 F.Y. conceived and led this research. J.S. performed the main research and managed technical implementa-
640 tion. J.S., F.Y., and J.S. (Junjie) designed the SA token and Saprot. Z.L. and J.S. developed the ColabSaprot
641 and SaprotHub. F.Y., J.S., and S.O. designed ColabSaprot and SaprotHub. S.O. and M.S. engaged in discus-
642 sions of ColabSaprot and SaprotHub, provided many constructive suggestions, and revised the manuscript.
643 C.H. conducted partial zero-shot mutation experiments. Y.Z. collected the AFDB dataset. Y.H., X.C. and
644 X.Z. conducted experiments on eTDG. D.M. participated in early wet lab experiments. J.S. and F.Y. wrote
645 the manuscript. OPMC authors contributed as described in the paper, including manuscript proofreading,
646 idea discussions, expert advice, promoting, testing and development of ColabSaprot.

647 Supplementary information

648 OPMC

649 The Open Protein Modeling Consortium (OPMC) is a collaborative initiative designed to unify the efforts
650 of the protein research community. Its mission is to facilitate the sharing and co-construction of resources,
651 with a particular focus on individually trained decentralized models, thereby advancing protein modeling
652 through collective contributions. OPMC offers a platform that supports a wide range of protein function
653 predictions, aiming to make advanced protein modeling accessible to researchers irrespective of their machine
654 learning expertise.

655 Here are some Q&A related to OPMC that we have gathered from our steering committee members and
656 other researchers: <https://theopmc.github.io/>

657 Be an OPMC member

658 The complete list of OPMC authors will be finalized in the last version of the paper. OPMC encompasses
659 two membership categories: the steering committee and regular members.

660 **Steering Committee:** The steering committee comprises esteemed researchers with doctoral degrees
661 who assume advisory roles within OPMC. Their primary responsibilities include:

- 662 • Providing constructive suggestions for OPMC and Saprothub.
- 663 • Supporting academic activities such as organizing workshops, tutorials, and participating in peer review
664 processes.
- 665 • Encouraging direct involvement of team members in the development of Saprothub or OPMC.
- 666 • Participating in the review process for new members joining OPMC.
- 667 • Contributing high-quality open-source datasets and models to enhance OPMC's resources.
- 668 • Making direct contributions to this paper.
- 669 • Making other contributions or providing academic services that drive the progress of OPMC.

670 **Regular Members:** Regular members actively contribute to OPMC through practical means. This
671 category includes independent researchers and professionals who:

- 672 • Participate in research, development, and collaborative efforts to advance OPMC's goals and objectives.
- 673 • Create innovative adapters for Saprothub, trained on either private or publicly available data.
- 674 • Contributing more high-quality datasets to enrich the available resources.
- 675 • Aid in the development and administration of Saprothub or ColabSaprot.

676 **Membership Eligibility:** We enthusiastically welcome individuals who make substantial contributions
677 to OPMC and Saprothub. The eligibility of regular members will be determined by the steering committee.
678 If you are interested in becoming a member of OPMC, please reach out to us. Note that all OPMC members
679 will be included as authors in the final version of this paper. The final list of members is expected to be
680 finalized within the next 4 to 12 months.

Dataset	Evaluation Metric	Train	Valid	Test
Training with AF2 structure				
Thermostability [45]	Spearman's ρ	5310	706	706
Subcellular localization [43]	Acc	10414	1368	1368
Binary localization [43]	Acc	6707	698	807
Training with PDB structure				
Structural similarity [52]	Spearman's ρ	300699	4559	4850
Structure class [52]	ACC	7990	955	1005
Metal ion binding [50]	ACC	5797	719	719
Binding site detection [52]	Mcc	2368	442	464
Training without structure				
Fluorescence [58]	Spearman's ρ	20963	5235	25517
Stability [58]	Spearman's ρ	53614	2512	12851
AAV [45]	Spearman's ρ	22246	2462	50432
β-lactamase [66]	Spearman's ρ	4158	520	520

Supplementary Table 1: Downstream dataset descriptions after all data pre-processing. All datasets that require structures as input are split based on structure similarity with 70 as threshold. For **Training with AF2 structure**, proteins with structures not available in AlphaFoldDB were excluded from both training and testing evaluations.

Model	ProteinGym	Mega-scale	ClinVar
	Spearman's ρ	Spearman's ρ	AUC
ESM-2 (650M) [25] (Science, 2023)	0.414	0.478	0.862
ESM-2 (15B) [25] (Science, 2023)	0.401	0.467	0.843
ProtBert [6] (TPAMI, 2021)	0.350	0.341	0.862
ESM-1v [54] (NeurIPS, 2021)	0.374	0.374	0.891
Tranception L [55] (ICML, 2022)	0.374	0.307	0.845
MSA Transformer [59] (ICML, 2021)	0.421	0.423	0.854
EVE [8] (Nature, 2021)	0.433	0.260	0.878
Saprot	0.457	0.574	0.909

Supplementary Table 2: More baseline results on zero-shot mutational effect prediction.

Task	Category	Metric	ESM-2	ProtBert	MIF-ST [67]	GearNet[69]	Saprot
Thermostability		Spearman's ρ	0.673	0.621	0.678	0.454	0.697
Subcellular localization	Training w/ AF2 structure	Acc%	84.67	78.65	82.46	75.73	86.18
Binary localization		Acc%	93.56	93.31	92.08	90.72	94.67
Structure class		Acc%	79.90	71.57	71.63	74.52	80.80
Structural similarity	Training w/ PDB structure	Spearman's ρ	0.677	0.601	0.555	0.101	0.695
Binding site detection		Mcc	0.540	0.459	0.504	0.415	0.602
Metal Ion binding		Acc%	74.97	68.19	69.58	67.61	79.97
Fluorescence		Spearman's ρ	0.698	0.698	-	-	0.699
Stability	Training w/o structure	Spearman's ρ	0.780	0.753	-	-	0.779
β -lactamase		Spearman's ρ	0.904	0.881	-	-	0.890
AAV		Spearman's ρ	0.813	0.783	-	-	0.829

Supplementary Table 3: More baseline results on supervised learning tasks. MIF-ST and GearNet are not applicable when protein structures are not available as input.

Task / Dataset	Metric	ProstT5	Saprot
ClinVar	AUC	0.620	0.909
Mega-scale	Spearman's ρ	0.194	0.574
Subcellular localization	Acc%	83.26	86.18
Binary localization	Acc%	94.11	94.67
Protein inverse folding	pLDDT	76.24	77.96
	RMSD	2.28	2.07
	TM-score	0.78	0.82
	Sequence recovery	33.72	36.15

Supplementary Table 4: Comparison between ProstT5 and Saprot in several tasks. ProstT5 is a seq2seq T5 model where the input, if it is an amino acid sequence, corresponds to the output of the 3Di sequence, and vice versa. In general, Saprot is comparable to ProstT5 in the protein inverse folding task, but substantially outperform it on other tasks, e.g., zero-shot mutational effect prediction and supervised learning tasks.

Supported Tasks	
Protein-level Classification	Active Site Prediction Post-Translational Modification (PTM) Site Prediction Mutation Impact Prediction Transmembrane Region Prediction Conservation Score Prediction Interface Residue Prediction Enzyme Binding Site Prediction Metal Binding Site Prediction Disease-Associated Variant Prediction Immune Receptor Binding Prediction Glycosylation Site Prediction
Post-translational Modification Site Classification	
Fold Class Prediction	
Localization Prediction	
Subcellular Localization	
Domain Classification	
Membrane Protein Classification	
Enzyme Classification	
Protein Stability Classification	
Protein Solubility Classification	
Antigenicity Classification	
Protein Aggregation Classification	
Protein Family Classification	
Protein Structural Class Prediction	
Protein Subunit Interface Classification	
Protein Motif Classification	
Protein Epitope Classification	
Protein Interaction Network Classification	
Protein Toxicity Classification	
Protein Expression Level Classification	
Protein Dynamics Classification	
Protein Flexibility Classification	
Protein Function Annotation	
Protein Interaction Type Classification	
Protein Disease Association Classification	
Protein Functional Domain Classification	
Protein-protein Classification	Protein-Protein Interaction (PPI) Prediction Interaction Type Classification Disease-Associated Interaction Prediction Mutational Impact Classification Functional Annotation Classification Interaction Network Module Detection Interaction Stability Prediction
Protein-protein Regression	Interaction Strength Prediction Binding Free Energy Calculation Interaction Affinity Prediction Mutational Impact Quantification Interaction Kinetics Prediction
Zero-shot mutational Effect prediction	Enzyme Activity Prediction Virus Fitness Prediction Driver Mutation Prediction Protein Stability Prediction Immune Escape Prediction MHC Binding Prediction Pathogenicity Prediction Membrane Protein Function Prediction Functional Domain Impact Prediction
Protein-level Regression	Thermal Stability Prediction Fluorescence Intensity Prediction Binding Affinity Prediction Solubility Prediction Expression Level Prediction Enzymatic Activity Prediction Hydrophobicity Prediction Protease Sensitivity Prediction Protein Interaction Strength Prediction Drug Binding Affinity Prediction Protein Stability Change upon Mutation Prediction Protein Ligand Binding Kinetics Prediction Protein Stability in Different Temperatures Prediction Protein Stability in Different Solvents Prediction
Residue-level Classification	Secondary Structure Prediction Binding Site Prediction
	Enzyme Function Optimization Protein Stability Enhancement Protein Folding Prediction Antibody Design Development of Novel Protein Drugs Vaccine Design Protein Engineering Synthetic Antibody Fragments Design

Supplementary Table 5: Tasks that Saprot potentially supports. As a near-universal protein language model, Saprot potentially supports hundreds of protein prediction tasks (also see Fig. 1 c,d). These tasks include, but are not limited to: protein-level classification, protein-level regression, residue-level classification, protein-protein classification, protein-protein regression, zero-shot mutational effect prediction, and protein sequence design.

Mutation	Predicted score	T>G% (wildtype=10.3865)	Enhancement factor
I37P	4.8993	15.7850	1.5198
L74A	4.7470	12.9784	1.2495
R8M	3.8026	14.7522	1.4203
L74E	3.6393	20.4948	1.9732
S5M	3.6230	15.0282	1.4469
P43K	3.3585	14.7687	1.4219
P6M	3.3499	15.5822	1.5002
P43R	3.2832	13.6750	1.3166
I37S	3.2649	16.4798	1.5866
I103E	3.2008	11.9370	1.1493
H11P	3.1735	14.6748	1.4129
H11K	3.1643	18.0260	1.7355
I37R	3.1442	14.8948	1.4341
I103D	3.1406	6.9425	0.6684
F1M	3.0996	6.0477	0.5823
L74Q	3.0893	19.1058	1.8395
H11A	2.9147	17.1734	1.6534
P4M	2.7864	15.2689	1.4701
I37A	2.7000	16.3456	1.5737
K252L	2.6721	1.2497	0.1203

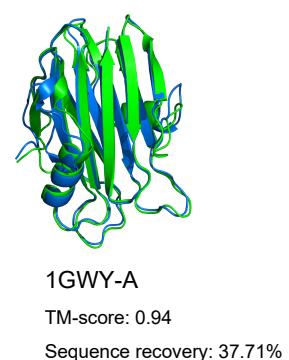
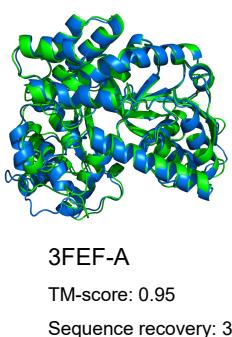
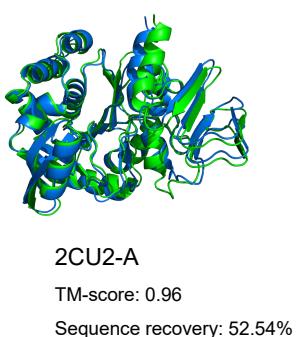
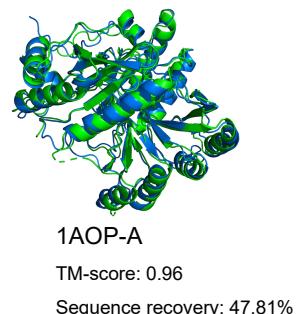
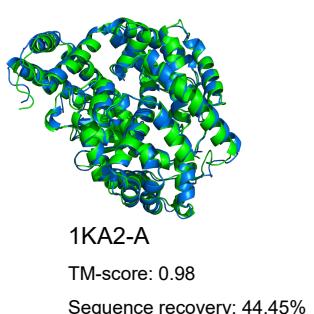
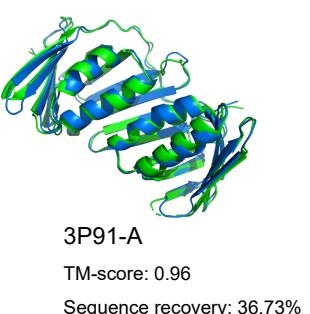
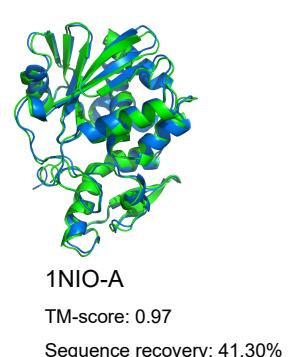
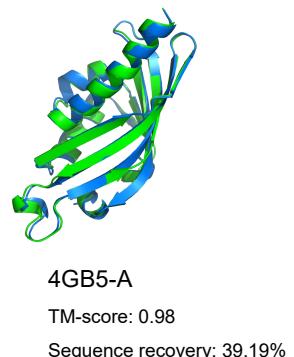
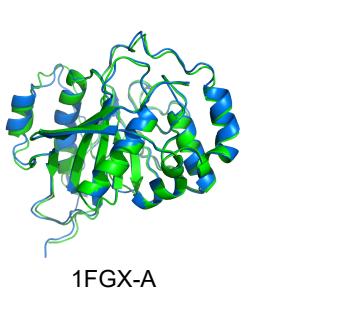
Supplementary Table 6: The editing efficiency of T-to-G mutations in the top 20 variants predicted by ColabSaprot. Among the top 20 variants, 17 mutations exhibited higher editing efficiency compared to the wild-type. Specifically, the editing efficiency at the L47E, L74Q and H11K sites were nearly doubled compared to the wild-type.

Benchmark Scores

<https://proteingym.org/benchmarks>

Rank▲	Model name	Model type	Avg. Spearman	Rank▲	Model name	Model type	Avg. Spearman	Rank▲	Model name	Model type	Avg. Spearman	Rank▲	Model name	Model type	Avg. Spearman
1	SaProt (650M)	Protein language model	0.457	11	ProtSSN (k=30 h=512)	Hybrid - Structure & PLM	0.438	21	ProtSSN (k=10 h=768)	Hybrid - Structure & PLM	0.423	31	ESM2 (15B)	Protein language model	0.401
2	TranceptEVE L	Hybrid - Alignment & PLM	0.456	12	ProtSSN (k=30 h=768)	Hybrid - Structure & PLM	0.438	22	ESM-IF1	Inverse folding model	0.422	32	MIF-ST	Hybrid - Structure & PLM	0.4
3	TranceptEVE M	Hybrid - Alignment & PLM	0.455	13	ProtSSN (k=30 h=1280)	Hybrid - Structure & PLM	0.438	23	MSA Transformer (single)	Hybrid - Alignment & PLM	0.421	33	EVmutation	Alignment-based model	0.395
4	GEMME	Alignment-based model	0.455	14	VESPA	Protein language model	0.436	24	DeepSequence (ensemble)	Alignment-based model	0.419	34	ESM-1b	Protein language model	0.394
5	TranceptEVE S	Hybrid - Alignment & PLM	0.452	15	Tranception L	Hybrid - Alignment & PLM	0.434	25	Tranception S	Hybrid - Alignment & PLM	0.418	35	VESPAI	Protein language model	0.394
6	ProtSSN (ensemble)	Hybrid - Structure & PLM	0.449	16	MSA Transformer (ensemble)	Hybrid - Alignment & PLM	0.434	26	ESM2 (650M)	Protein language model	0.414	36	Progen2 XL	Protein language model	0.391
7	ProtSSN (k=20 h=1280)	Hybrid - Structure & PLM	0.442	17	ProtSSN (k=10 h=1280)	Hybrid - Structure & PLM	0.433	27	DeepSequence (single)	Alignment-based model	0.407	37	ESM2 (150M)	Protein language model	0.387
8	ProtSSN (k=20 h=512)	Hybrid - Structure & PLM	0.441	18	EVE (single)	Alignment-based model	0.433	28	ESM-1v (ensemble)	Protein language model	0.407	38	MIF	Inverse folding model	0.383
9	ProtSSN (k=20 h=768)	Hybrid - Structure & PLM	0.44	19	ProtSSN (k=10 h=512)	Hybrid - Structure & PLM	0.43	29	SaProt (35M)	Protein language model	0.406	39	Progen2 L	Protein language model	0.38
10	EVE (ensemble)	Alignment-based model	0.439	20	Tranception M	Hybrid - Alignment & PLM	0.427	30	ESM2 (3B)	Protein language model	0.406	40	Progen2 M	Protein language model	0.379

Supplementary Fig. 1: Saprot (checkpoint: https://huggingface.co/westlake-repl/SaProt_650M_AF2) ranked first on the public ProteinGym leaderboard until May 2024, surpassing over 60 well-known methods. Note that Saprot is a single model, while most other top-ranked models are hybrid models.



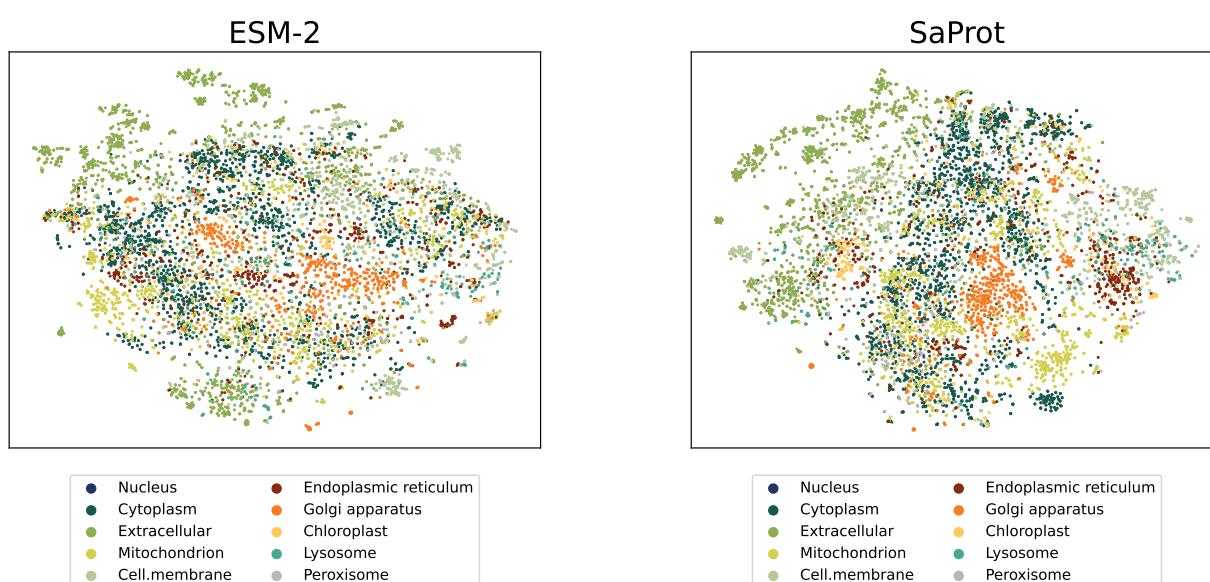
Supplementary Fig. 2: Case studies of protein inverse folding using Sapro. The blue-colored protein structures are predicted using ESMFold, with residue sequences generated by Sapro, and the green-colored structures are the experimentally determined structures. Given the backbone, Sapro is capable of predicting diverse sequences with few residue overlap while maintaining almost the same 3D structure.



(a) Embedding visualizations of ESM-2 650M version and Saprot 650M version on SCOPe database. We use the non-redundant version ($PIDE < 40\%$) of the SCOPe database.

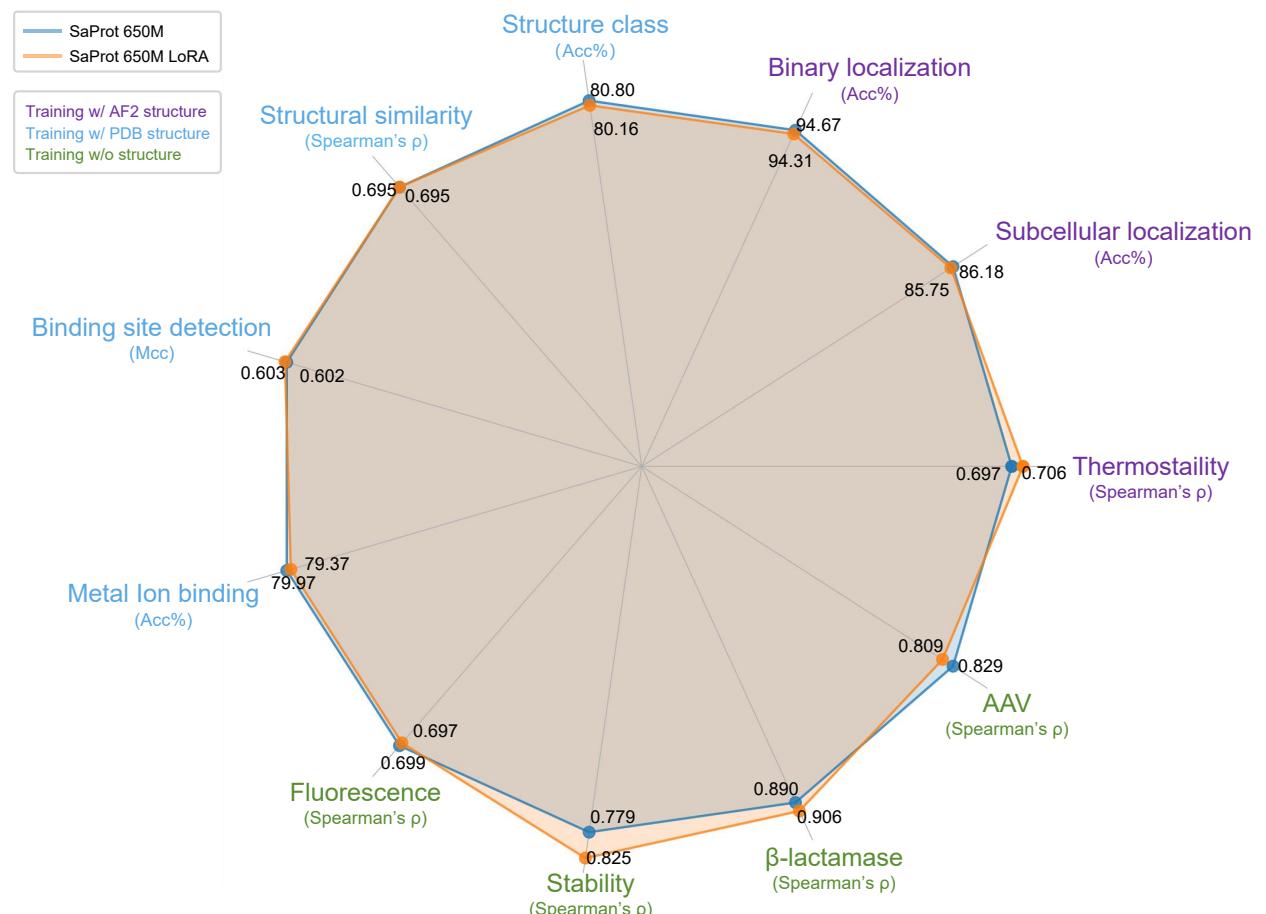


(b) Embedding visualizations of ESM-2 and Saprot on binary localization dataset.



(c) Embedding visualizations of ESM-2 and Saprot on subcellular localization dataset.

Supplementary Fig. 3: Embedding visualization



Supplementary Fig. 4: Comparison between full model fine-tuning and parameter efficient fine-tuning. While greatly reducing trainable parameters, fine-tuning Saprot using adapter module achieves competitive performance on all experimental settings, which demonstrates its applicability on various biological tasks.